

# **The Performance Evaluation Research Center (PERC)**

## **Participating Institutions:**

Argonne Natl. Lab.

Univ. of California, San Diego

Lawrence Berkeley Natl. Lab.

Univ. of Illinois

Lawrence Livermore Natl. Lab.

Univ. of Maryland

Oak Ridge Natl. Lab.

Univ. of Tennessee, Knoxville

**Website: <http://perc.nersc.gov>**



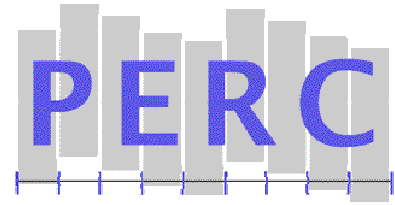
# PERC Overview

- An “Integrated Software Infrastructure Center” (ISIC) sponsored under DoE’s SciDAC program.
- Funding: approx. \$2.4 million per year.
- Mission:
  - Develop a science of performance.
  - Engineer tools for performance analysis and optimization.
- Focus:
  - Large, grand-challenge calculations, especially large-scale scientific codes used in SciDAC projects.



# Specific Objectives

- Understand key factors in scientific codes that affect performance.
- Understand key factors in computer systems that affect performance.
- Develop models that accurately predict performance of codes on systems.
- Develop an enabling infrastructure of tools for performance monitoring, modeling and optimization.
- Validate these ideas and infrastructure via close collaboration with DOE Office of Science and others.
- Transfer the technology to end users.



# Economic Benefits

Consider the economic value of improving the performance of a single high-end scientific application code by 20%.

Assume:

- \$10 million computer system lease cost per year.
- \$10 million per year in site costs, support staff, etc.
- 10-year lifetime of code.
- Code uses 5% of system cycles each year.

Savings: \$2,000,000.

Scientific benefit (additional computer runs and research) is probably much higher.



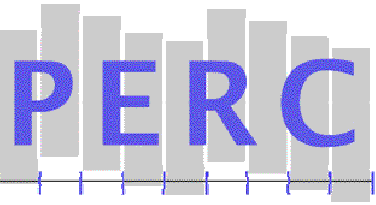
# Quantitative Feedback to Vendors

- We are invited by vendors to provide guidance on the design of current and future systems.

BUT

- At present we can provide only vague information – little if any quantitative data or rigorous analysis.

The performance monitoring and modeling capability being developed in PERC will significantly improve our ability to influence future scientific computer systems.



# Four Thrusts

## **Better Benchmarks:**

- Kernel benchmarks extracted from real codes reduce complexity of analyzing full-size benchmarks.
- Low-level benchmarks measure key rates of data access at various levels of memory hierarchy.

## **Modern performance monitoring tools:**

- Flexible instrumentation systems capture hardware and software interactions, instruction execution frequencies, memory reference behavior, and execution overheads.
- An advanced data management infrastructure tracks performance experiments and data across time and space.

## Performance modeling:

- *Application signature* tools characterize applications independent of the machine where they execute.
- *Machine signature* tools characterize computer systems, independent of the applications.
- *Convolution* tools combine application and machine signatures to provide accurate performance models.
- *Statistical models* find approximate performance models based on easily measured performance data.
- *Performance bound* tools determine ultimate potential of an application on a given system.

## Performance optimization:

- *Compile-time optimization* mechanisms analyze source code to improve performance.
- *Self-tuning software* automatically tunes code based on real-time measurements of hardware environment.
- *Performance assertions* permit user-specified run-time tests to possibly change the course of the computation depending on results.
- *Performance portability* programming techniques to insure that code runs at near-optimal performance across a variety of modern systems.

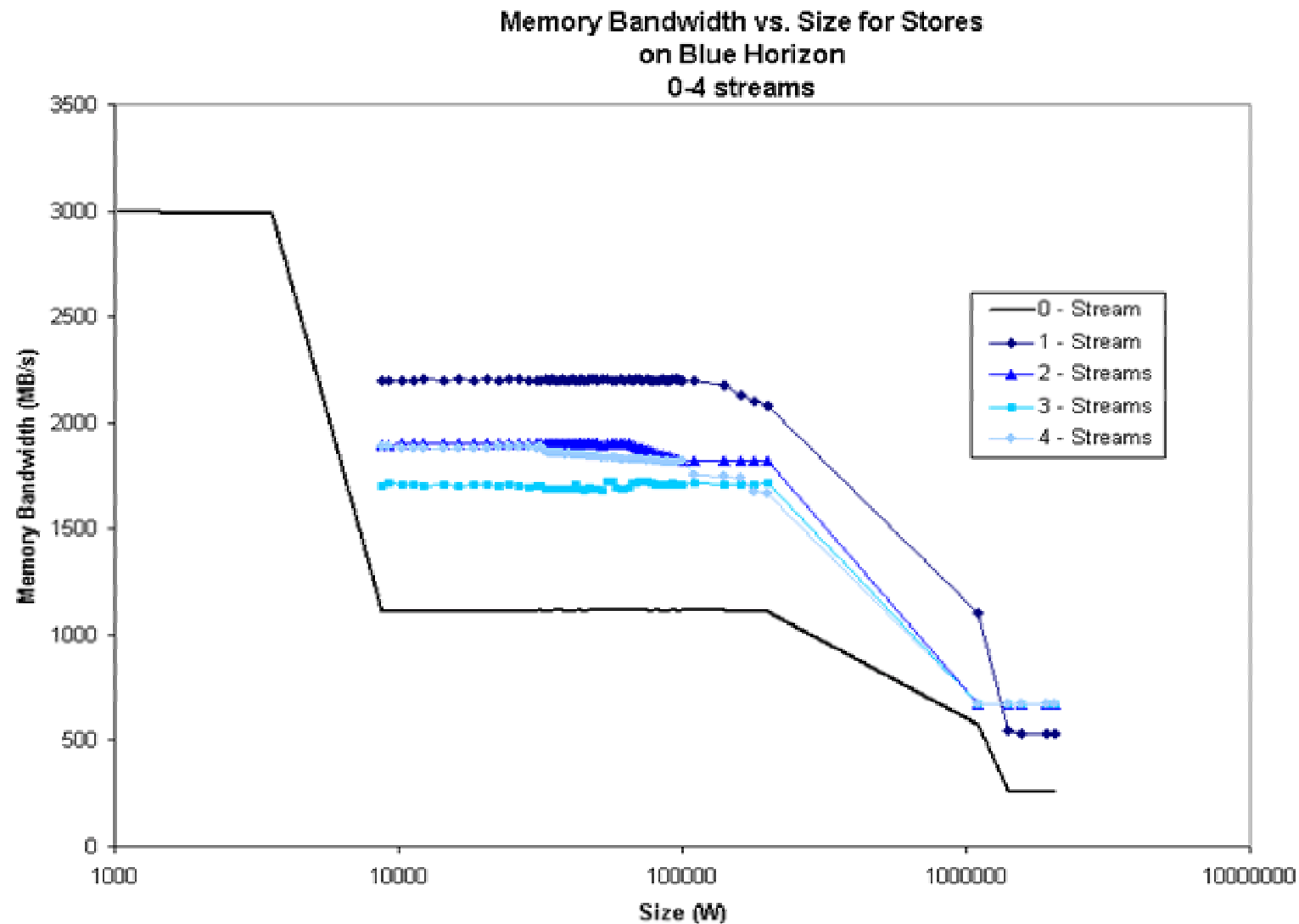




# Partnerships with SciDAC Scientific Projects

- Terascale Simulation of Neutrino-Driven Supernovae
- Advanced Computing for 21st Century Accelerators
- National Computational Infrastructure for Lattice Gauge Theory
- Collaborative Design and Development of Community Climate System Model for Terascale Computers
- Numerical Computation of Wave-Plasma Interactions
- The Plasma Microturbulence Project
- Terascale Optimal PDE Solvers

# Memory Access Patterns (MAPS)





# Performance Model Vs Actual

# CPUs	Real Time	Predicted Time	% Error
2	31.78	31.82	0.13
4	29.07	31.27	7.57
8	36.13	33.72	6.67
64	44.91	43.91	2.23
96	48.87	47.15	3.52
128	52.88	52.46	0.79

PETSc kernel, run on IBM SP at SDSC.



# SvPablo Graphical Interface

**svPablo**

Project Instrument View GenCallGraph Help

**Project Description:** PCTM on IBM-SP (seaborg)

**Source Files:** fod.F  
fod\_setup.F

**Performance Contexts:**  
IBM-SP, 64 procs, other Metrics  
IBM-SP, 64 procs, other Metrics, 10 days  
IBM-SP, 32 procs, other Metrics, 10 days  
IBM-SP, 16 procs, other Metrics, 10 days  
IBM-SP, 8 procs, other Metrics, 10 days  
IBM-SP, 4 procs, other Metrics, 10 days

**Routines in Source File:** fod  
fod\_setup  
fod\_timer\_clear  
fod\_timer\_start  
fod\_init

**Routines in Performance Data:** ccm3  
oceanstep  
icesten  
mpi\_c  
mpi\_c

**Source File:** /u0/cmendes/PCTM/pom2/src/sources/fod.F

**Specific Metric**  
Call Statistics count:  
240.0000 -- ccm3  
Dismiss Help

**Specific Metric**  
Call Statistics Duration:  
211.2399 -- ccm3  
Dismiss Help

**Specific Metric**  
HW Statistics by Line Floating Point Instructions:  
12295122047.0000 -- ccm3  
Dismiss Help

**Specific Metric**  
HW Statistics by Line Load Misses in D1:  
300348320.0000 -- ccm3  
Dismiss Help

**Specific Metric**  
HW Statistics by Line Branch Instructions:  
6944481284.0000 -- ccm3  
Dismiss Help

**Specific Metric**  
HW Statistics by Line TLB misses:  
151118598.0000 -- ccm3  
Dismiss Help

**Legend: Source Code Metrics**

Column 1: Call Statistics count  
240  
10

Column 2: Call Statistics Duration  
211.24  
119.525

Column 3: Loop Statistics count  
0  
0

Column 4: Loop Statistics Duration  
0  
0

Column 5: HW Statistics by Line Floating Point Instr  
1.65722e+10  
0

Column 6: HW Statistics by Line Load Misses in D1  
8.62196e+08  
2.46132e+08

Column 7: HW Statistics by Line Branch Instructions  
6.94448e+09  
0

Column 8: HW Statistics by Line Load Instructions  
3.10653e+10  
0

Column 9: HW Statistics by Line Instruction Cache M  
9.55654e+07  
6.80955e+07

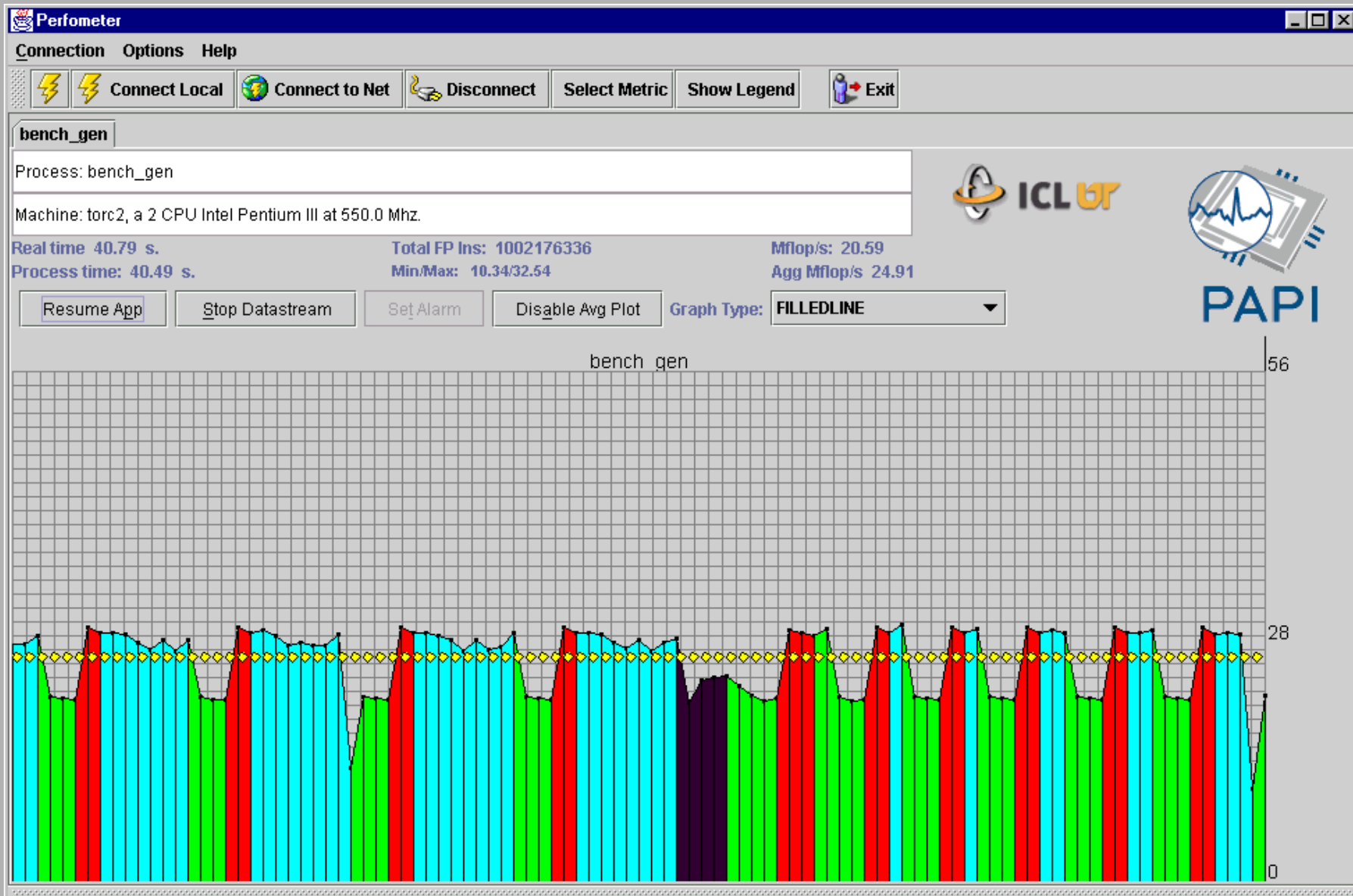
Column 10: HW Statistics by Line TLB misses  
1.51119e+08  
1.0256e+07

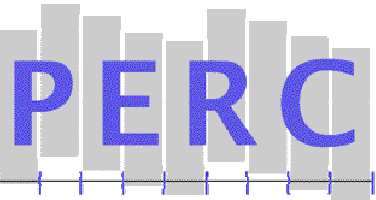
Instrument/Clear Line View Line Data

Dismiss Help



# PAPI Perfometer Interface



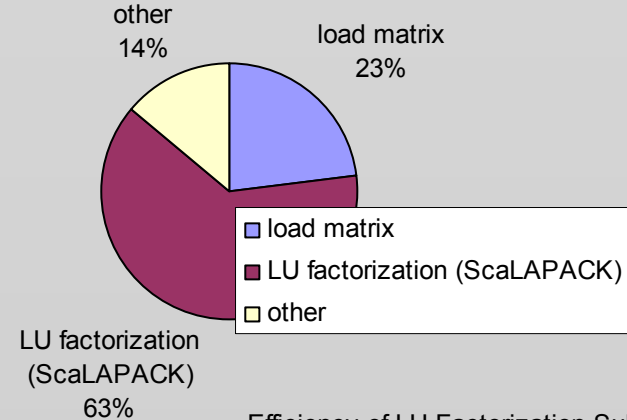


# Performance Analysis

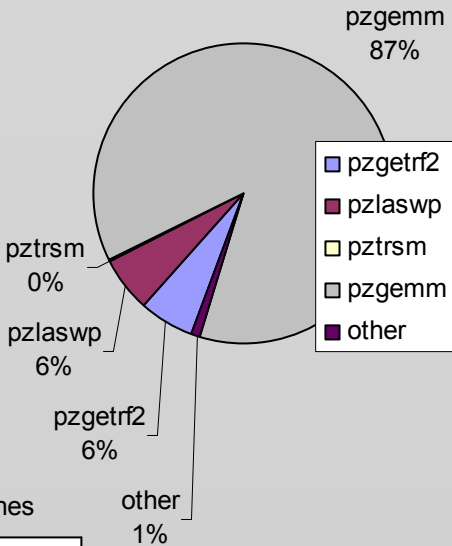
## AORSA3D (fusion)

AORSA3D was ported and benchmarked on IBM and Compaq platforms. A detailed performance analysis has begun using SvPablo and PAPI. The results below are for a 400 Fourier mode run on 16 processors and 1 node of an IBM SP (Nighthawk II / 375MHz).

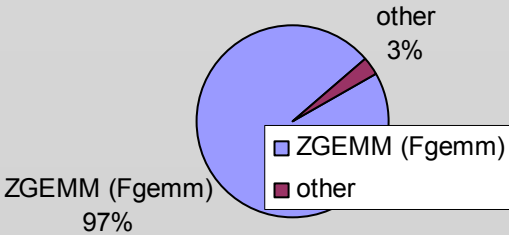
Time Profile of Total Execution



Time Profile of LU Factorization



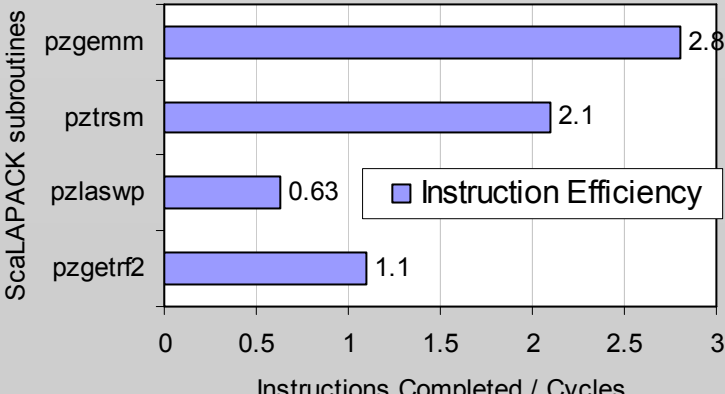
Time Profile of PZGEMM



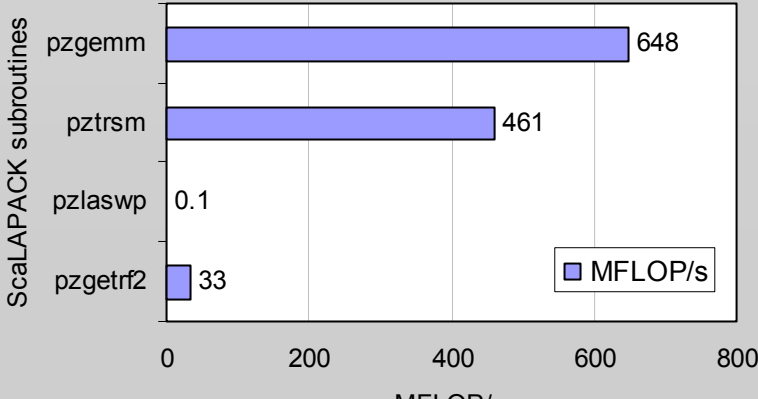
### Performance for ZGemm

Denisty of Mem Access	1
Denisty of FLOPs	1.8
MFLOP/s	664
L1 cache hit rate	0.98
L2 cache hit rate	0.96
TLB misses	285034

Efficiency of LU Factorization Subroutines



MFLOP Rates for LU Factorization Subroutines

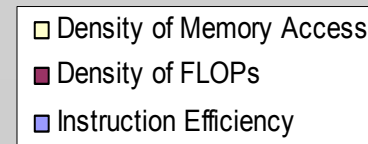
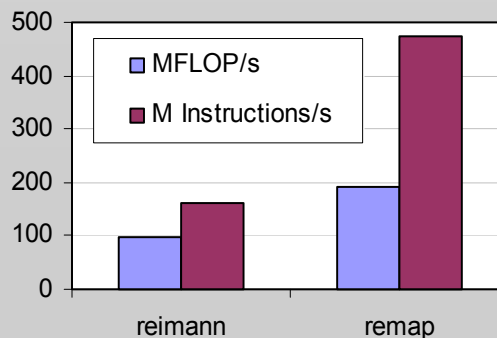
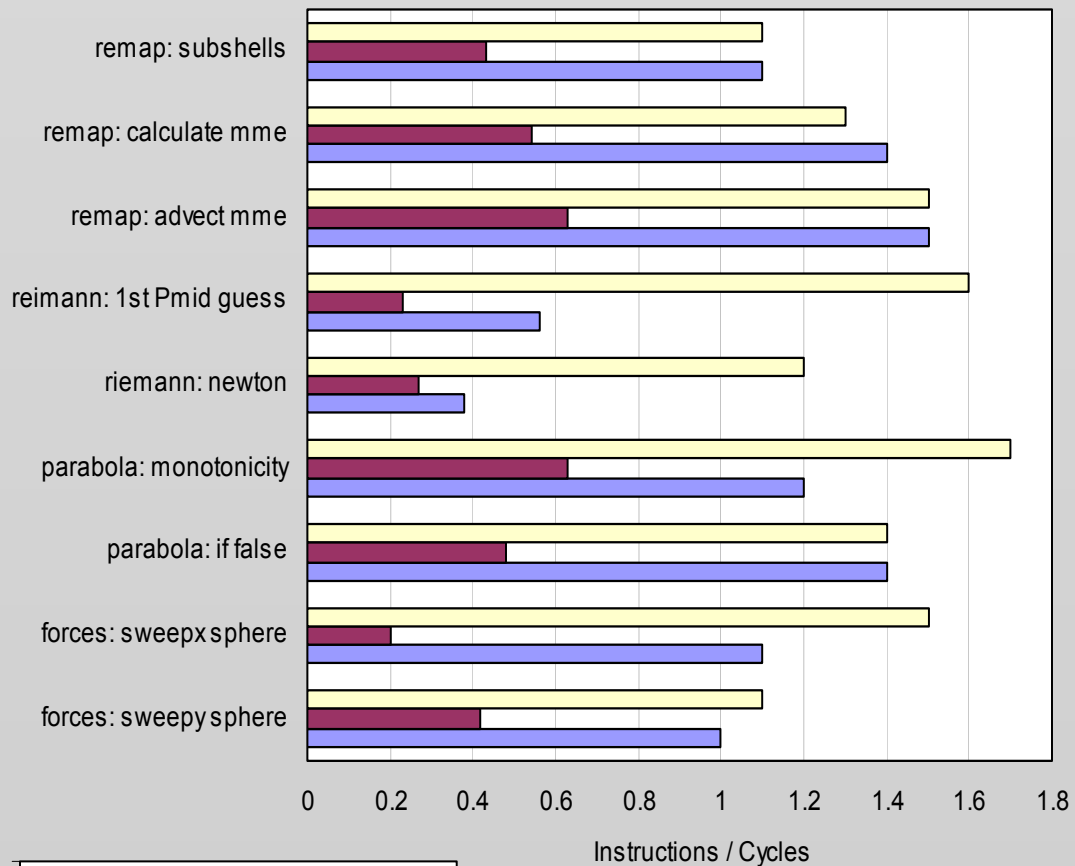
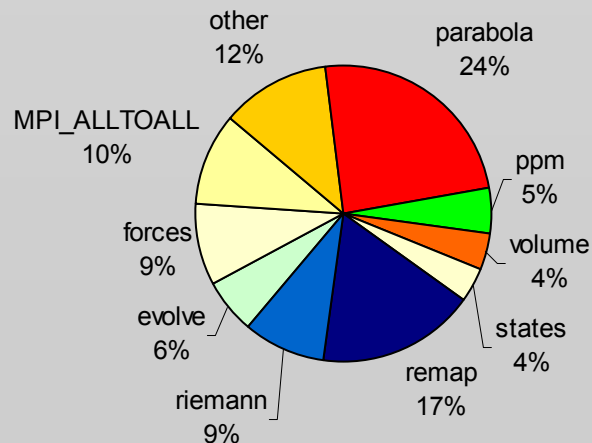
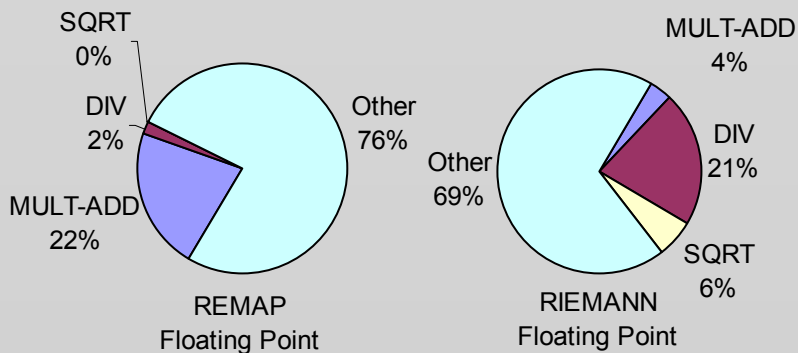


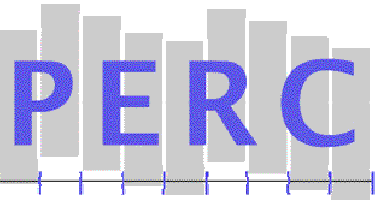


# Performance Analysis

## EVH1 (high-energy physics)

Aggregate performance measures over all tasks for a .1 simulation-second run. Collected with PAPI on an IBM SP (Nighthawk II / 375MHz).

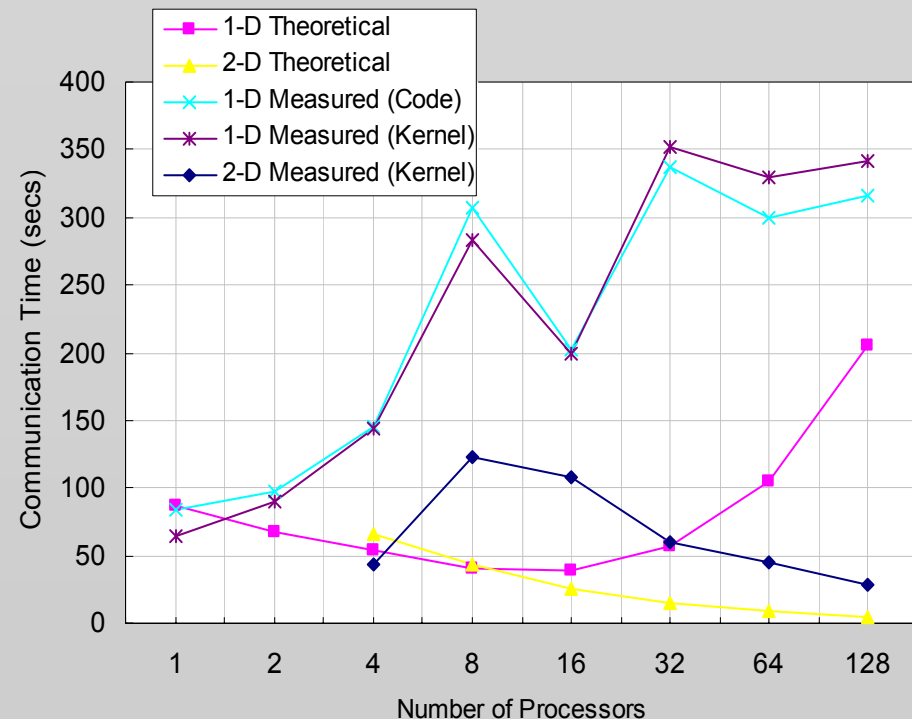
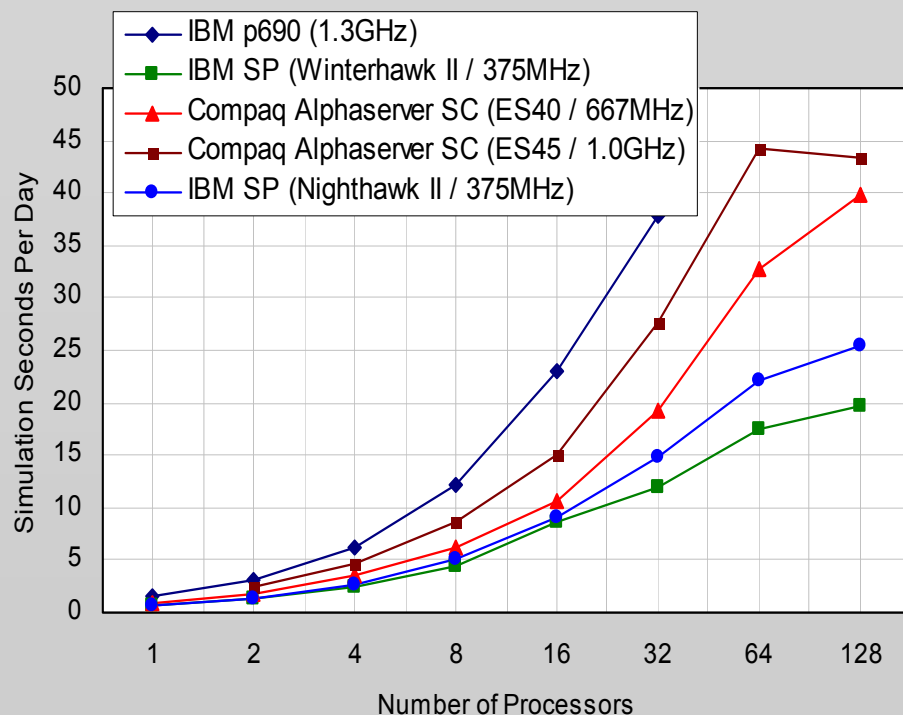




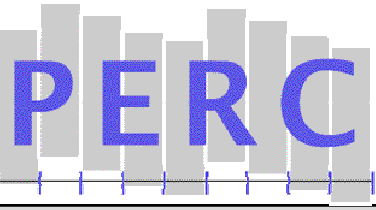
# Analysis and Tuning EVH1 (high-energy physics)

It was found that for more than 64 processors, EVH1 is communication-bound. The predominant routine (>50% of execution time) at this scale is MPI\_ALLTOALL, which is used in matrix-transpose-like operations.

The current implementation uses a 1D decomposition for the matrix-array; a modeling and analysis study has shown that a 2D decomposition would result in a large improvement. Benchmarking and analysis results are below.





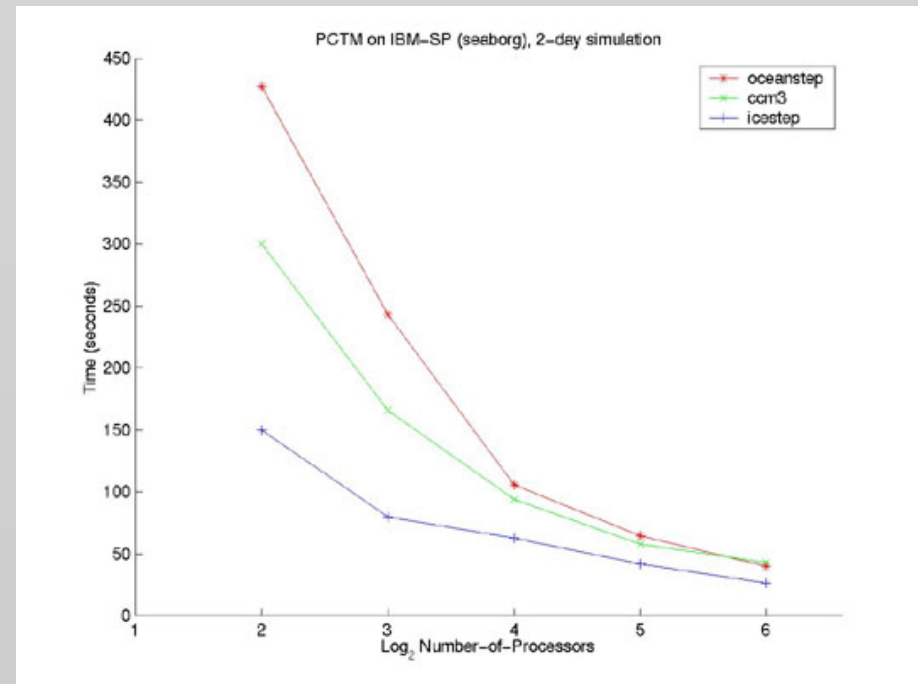
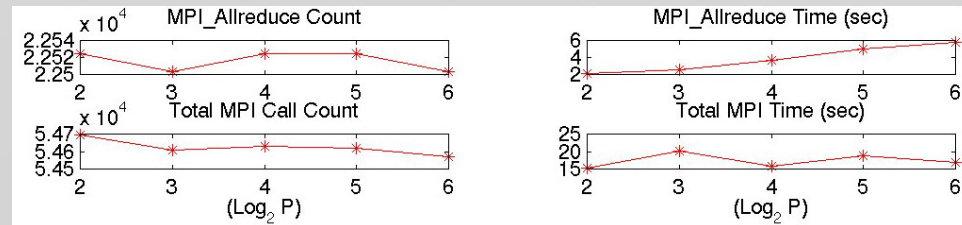
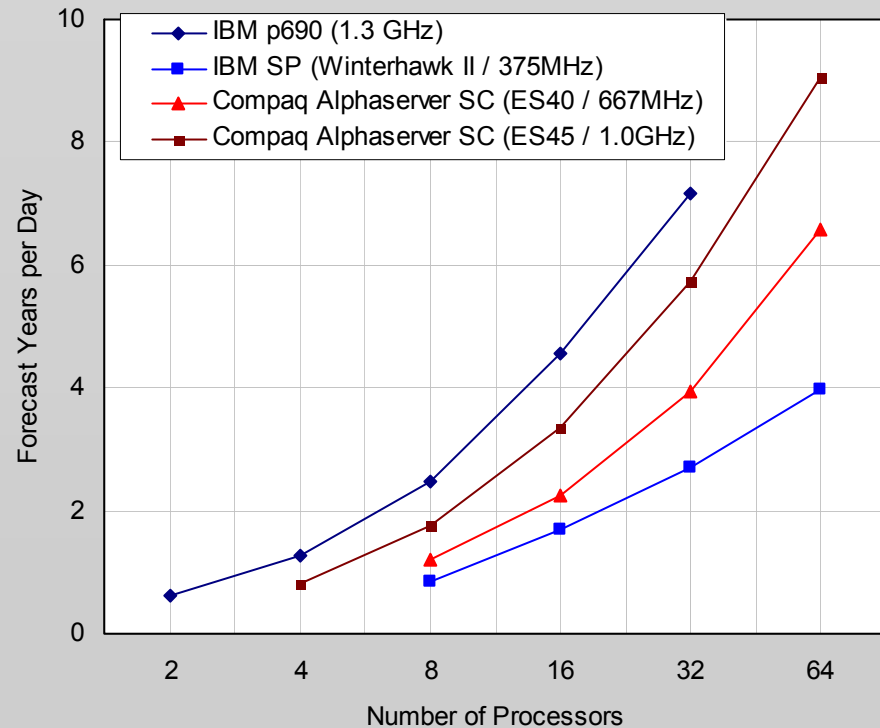


# Performance Analysis

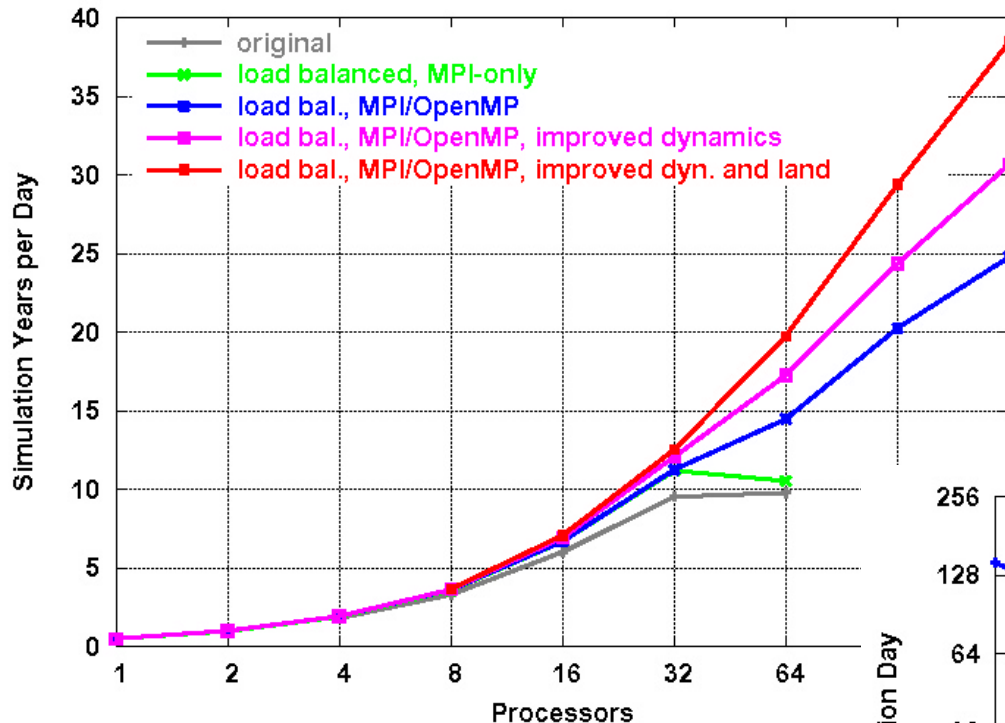
## PCTM (climate modeling)

In studies with PCTM, we examined the impact of different domain decompositions, as well as using fewer MPI processes per SMP node. These studies indicated a strong performance dependence on message-passing performance, where using 64 MPI processes on 32 4-way SMP nodes (leaving 64 processors idle) was 30% faster than when running on 16 SMP nodes.

PCTM Simulation Years per Day

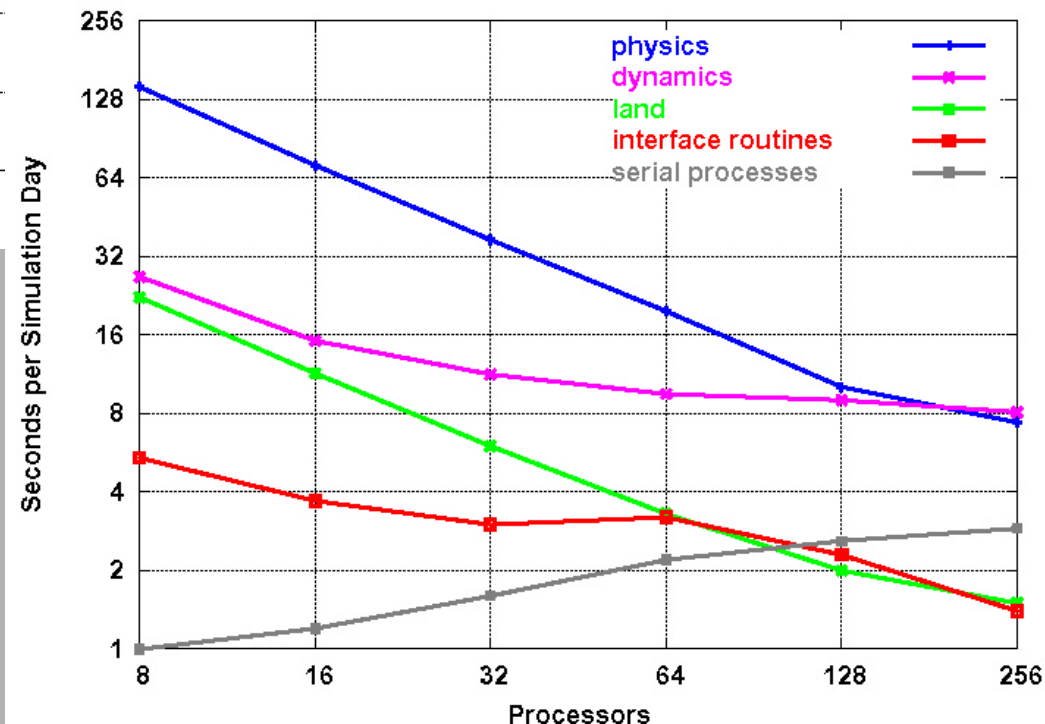


# Analysis and Tuning CAM (atmospheric modeling)



CAM performance measurements on IBM p690 cluster (and other platforms) were used to direct development process. Graph shows performance improvement from performance tuning and recent code modifications.

Profile of current version of CAM indicates that improving the serial performance of the physics is the most important optimization for small numbers of processors, and introducing a 2D decomposition of the dynamics (to improve scalability) is the most important optimization for large numbers of processors.





# Working with PERC

- Benchmarking
  - Application group works with PERC to specify relevant benchmark codes and problems.
  - PERC characterizes performance, generates performance models, and suggests optimizations.
- Performance Tools
  - PERC trains application developers to use tools.
  - Application group uses tools in development, providing feedback on functionality and future development

**For further information: <http://perc.nersc.gov>**